

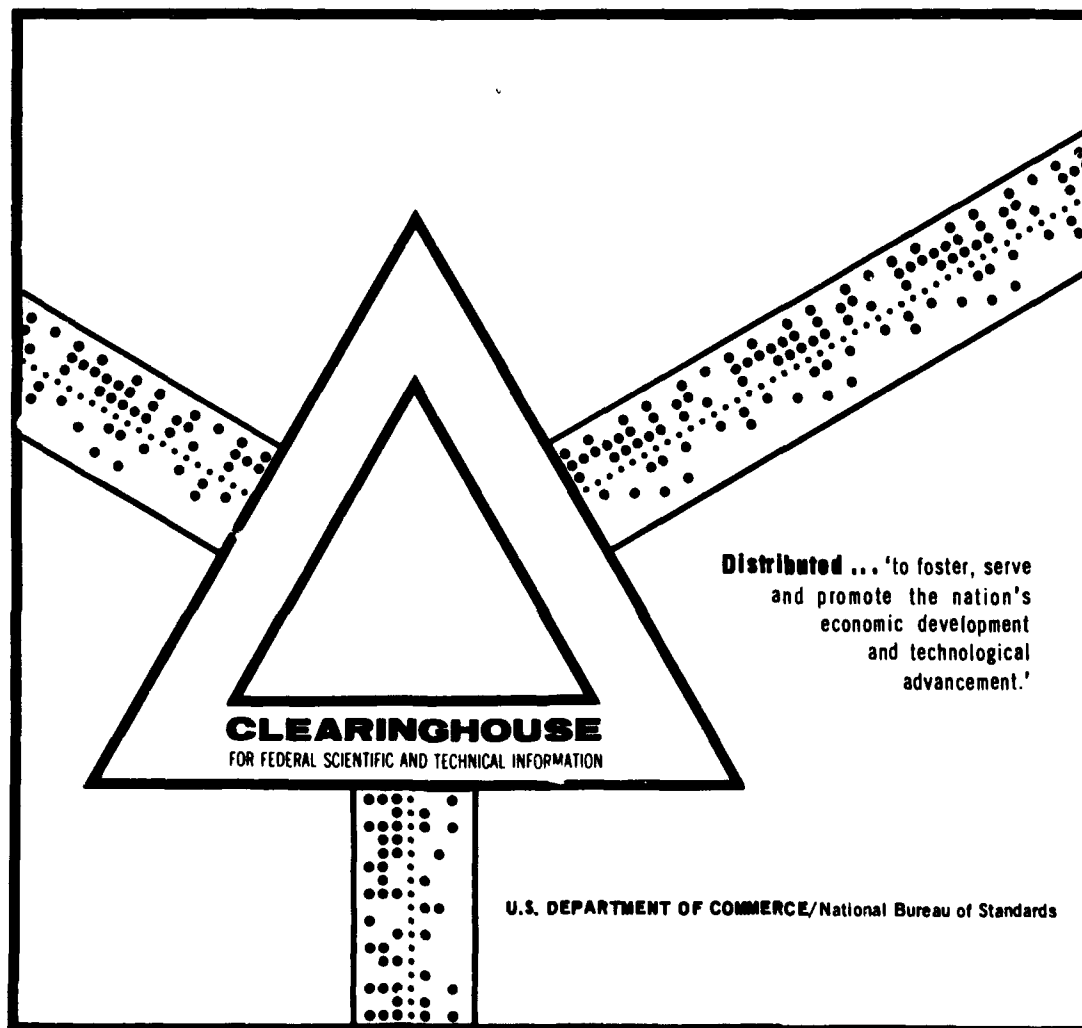
AD 700 924

**SOCIAL OPTIMIZATION VERSUS SELF-OPTIMIZATION
IN WAITING LINES**

I. Adler, et al

Technion - Israel Institute of Technology
Haifa, Israel

October 1969



This document has been approved for public release and sale.

AD700924

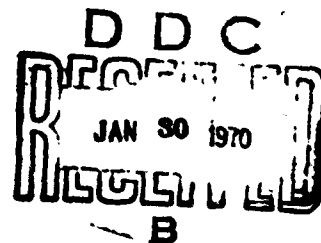
SOCIAL OPTIMIZATION VERSUS SELF-OPTIMIZATION
IN WAITING LINES

I. Adler and P. Naor

Operations Research, Statistics and Economics
Mimeograph Series No. 56

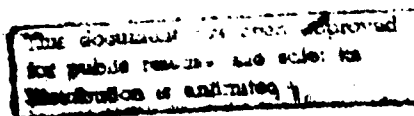


הטכניון — מכון טכנולוגי לישראל
הפקולטה להנדסת תעשייה וניהול



TECHNION — ISRAEL INSTITUTE OF TECHNOLOGY
FACULTY OF INDUSTRIAL AND MANAGEMENT ENGINEERING
HAIFA, ISRAEL

Reproduced by the
CLEARINGHOUSE
for Federal Scientific & Technical
Information Springfield Va. 22151



**SOCIAL OPTIMIZATION VERSUS SELF-OPTIMIZATION
IN WAITING LINES**

I. Adler and P. Naor

**Operations Research, Statistics and Economics
Mimeograph Series No. 56**

This research has been sponsored in part by the Logistics and Mathematical Statistics Branch, Office of Naval Research, Washington, D.C., under Contract F61052-68-C-0014.

This research was initiated in the Faculty of Industrial and Management Engineering, Technion-Israel Institute of Technology, Haifa, Israel and was completed in the Department of Operations Research, Stanford University, Stanford, Calif. The report is issued simultaneously at Stanford University.

October 1969

A B S T R A C T

A queueing model is considered where customers arriving in a Poisson stream are given the choice of either joining the waiting line or - by declining to do so - of foregoing the benefits accruing through service. The decision of each customer is based on his concrete benefit-cost analysis. Since his service time is constant, and exhaustive information as to the actual state of the system is available, both of the alternatives presented to the individual customer are completely deterministic and his decision is not reached under uncertainty or risk. The cost structure envisaged as well as additional assumptions give rise to a queueing model with limited waiting room apparently not previously considered in the literature. After detailed analysis of the model and blending with the cost structure it is shown that the criterion for self-optimization of the customers will not bring about social optimization, the latter being defined as the maximally feasible expected net gain per unit time accruing to the totality of customers. A number of simple and comprehensive optimization equations are derived. By marginal analysis the correctness of the simple equations is verified and their applicability is extended to models possessing more general character.

SOCIAL OPTIMIZATION VERSUS SELF-OPTIMIZATION
IN WAITING LINES

I. Adler and P. Naor

I. Introduction

The following problem has been discussed in a recent communication (Naor (1969)): Customers arrive in a Poisson stream at a service station; each of them is made aware of the current length of the queue, of the monetary reward he will attain through completion of service, and of his own waiting cost per unit time. Any customer is given the choice between two alternatives: joining the queue or balking the opportunity of doing so. It is desired to reach a rational decision between these alternatives. What is the proper criterion on which to base the decision? It was shown in that study that, if each customer reaches his conclusion on the basis of his narrow self-interest, a reasonable social objective function representing the public good will not be optimized. Rather it was shown that individuals acting solely on the basis of their self-interest will impose more congestion on the system than is socially desirable.

The particular model employed in the above study envisaged service time to be an exponentially distributed random variable. Hence, in such a situation the individual customer - on pursuing his self-interest - reaches his conclusion under risk; that is, he realizes that the random variable "future queueing time" a priori possesses an Erlang distribution (with known parameters) and in the comparison of outcomes of his possible decisions he makes use of expected queueing costs (among other items). Now in any actual realization it may turn

out, on hindsight, that the customer could have been better off, had he decided differently. One is tempted to consider the possibility that these a posteriori erroneous decisions are somehow associated with the contradiction between social optimization and self-optimization.

One purpose of the present investigation is to show that such is not the case. In this study a model is presented such that the individual customer is faced with a decision problem under certainly and yet his rational criterion of choice (under the specification of self-optimization) between the alternatives will not bring about social optimality. Again, as in the previous case, at some times the individual customer will be prepared to join the waiting line whereas overall considerations of optimality deem this prohibitive. It is the general rule that self-optimization tends to over-congest the system as compared with social optimization.

In the previous study algebra only was employed in the derivation of the various quantities of interest. The exponential distribution assumption with respect to the service time generated a situation where the state space was capable of description by a single integer only and the number of possible states was finite. This is advantageous if one's primary interest is to demonstrate the contradiction between self-optimization and social optimization. However, the case of exponential service distribution does not lend itself to convenient generalizations if one intends to investigate situations where the state space is infinite and comprises all real numbers in the interval $(0, v)$ ($v > 1$). In order to make progress toward such cases it is convenient to start with fixed and equal service times. The feasibility of attaining the higher degree of generality through the use of fixed (rather than exponentially distributed) service times contributes to the rationale underlying the approach presently taken in this communication. Finally, it is useful pointing out that the "pure" queueing problem (i.e. the stochastic model only, without the assumed cost structure) proposed here, to wit, the employment of a bounded and "non-integral" waiting room, seems not to have been dealt with in the literature. It is quite feasible that such models are of use in various applications other than those envisaged here. Thus, for instance, the mixing problem in chemical engineering stated and solved by Shinnar (1967) in queueing terms may be generalized on the lines of the pure stochastic model developed in the present paper.

II. Model Characteristics and Cost Structure

In the present Section we shall state the precise assumptions relating to the stochastic model as well as to cost structure. However, before going into the specifics it is useful to observe and state two distinctive qualities of the general setting in which the present model is situated.

1. Since the leitmotif of this study is contrasting two optimization procedures it is essential to define two distinct objective functions, one related to the aspirations of (non-cooperating) decision making individuals, the other concerned with the public good. There is no (mental) difficulty regarding the first position; the individual customer simply seeks to maximize his own net income. For a reasonable objective function to describe the second position one has to introduce a set of specific assumptions circumscribing the structure of public good. In the present communication we shall follow the mode employed in the preceding study (Naor (1969)) and choose the average net income accruing to the totality of customers in unit time as the objective function to be maximized. This choice of an objective function presupposes one of the following situations to prevail:
 - (i) There exists essentially only one genuine (overall) optimizer and individual customers are subcontractors of decision making, as it were, who are obliged by administrative fiat to reach their conclusions on the basis of rules prescribed by the overall optimizer.
 - (ii) Alternatively, a situation is envisaged where net gains of customers are considered to be comparative and additive, and by common agreement may be redistributed. Cooperation between customers - displayed by some through refraining from joining the queue apparently against their own best interests - will produce additional net income in unit time. This will be redistributed and, eventually, the average net profit accruing to each customer will exceed that of self-optimizing customers within a framework of non-cooperation. A feasible instrument

of control under the present circumstances is the imposition of a toll which, if wisely determined, will produce both social optimality and a reserve stock of money to be used for redistribution. Proper identification of a set of circumstances under which an overall objective function is deemed to exist is essential in an analysis of the present character.

2. In most queueing models with a built-in optimization procedure (e.g. through the agency of priority service rules or through control of the service intensity) it is assumed that all arriving customers - sooner or later - are going to be serviced. The feasible control actions (if such are envisaged at all) in most queueing models do not typically include the peremptory non-admission of a customer based on a cost-benefit analysis. While a number of models were developed which included the element of potential non-service - e.g. the balking and reneging models - this contingency was presented in probability terms only; non-admission was never considered to be an instrument of economic control. In the balking and reneging models the probability of not joining (and of leaving) the queue is associated with the customer's impatience - a psychological criterion rather than an economic one. Diverting the customer from the queue without rendition of any service is then a feasible course of action in the model area under consideration. To obviate difficulties - mostly of psychological origin - which stem from the feeling that customers must get some sort of service eventually, we may re-circumscribe the present framework in seemingly different terms: Customers may be served in two alternative, distinct modes. There is the standard mode of service which can always be relied upon and which is not associated with any queue of customers; it serves as reference point against which any other mode of service (concretely just one) may be compared. The non-standard mode of service is advantageous in monetary terms as compared with the standard mode if the waiting line of customers, ahead of the new arrival, is sufficiently small. If we describe the system in such terms - advantageous non-standard service with potential queue in front of the static.

compared with generally disadvantageous but queue-free standard service - we really deal with a model completely equivalent to the previous situation. To sum up, for a discussion of self-optimization versus social optimization to make sense non-admission of customers to the service station must be a feasible control action. Customers diverted from the station may be thought of either as not receiving service at all or as being rendered a standard type of service for which it is never necessary to queue up.

After these observations and the setting of the general framework the specifics of the stochastic model and of the cost structure may be stated in the following terms:

- (i) A stationary Poisson stream of customers - with parameter λ - arrives at a single service station.
- (ii) The service time necessary to satisfy and dispatch a customer is a constant T ; all service times are equal.
- (iii) On successful completion of service, the customer is endowed with a reward R (expressible in monetary units). All customer rewards are equal.
- (iv) The cost to a customer for staying in a queue (i.e. for queueing) is C monetary units in unit time. All customer costs are equal.
- (v) The newly arrived customer is required to choose one of two alternatives: either (a) he joins the queue, incurs the losses associated with spending some of his time in it, and finally obtains the reward; or (b) he refuses to join the queue - an action which does not bring about any gain or loss. The choice is made by the customer on comparing the net gains associated with each of these alternatives. Two modes of decision are examined: In one mode customers are assumed to act solely in their self-interest; it is sufficient for the net gain to the individual to be non-negative in order to induce him to join the queue. In the other decision mode each individual acts on behalf of the totality of customers and he assumes every customer to act in the same spirit; this customer seeks a decision criterion by which average net income in unit time is maximized.

Model assumption (i), (iii) and (iv) are identical with those appearing in the previous study. Assumption (v) is more inclusively formulated than its original counterpart in order to render equal status a-priori in the derivations to be carried out to both self-optimization and social optimization.

Assumption (ii) is radically different from its predecessor: Fixed - rather than exponentially distributed - service times are envisaged and this causes the individual customer to be placed in a deterministic decision situation. We note, furthermore, that this assumption (in conjunction with the others, of course) gives rise to a stochastic model which is interesting per se and may be put to use in other contexts as well. The mathematical techniques which have to be employed under the present assumption (ii) are of a different quality than those useful (and sufficient) if the assumption of exponentially distributed service times is considered valid. Finally, it will be shown that several of the results to be attained here by employing assumption (ii) serve as a more advantageous point of departure for some generalization than can be expected from the original model.

III. Finite and Nonintegral Queueing Capacity

What is the decision criterion of the individual customer seeking self-optimization? Clearly he reaches his decision under conditions of certainty. Upon his arrival he views the queue ahead of him which is made up of two parts: k customers are in the waiting line and one is in service. The outstanding service time of the latter is observed to be τ ($0 \leq \tau \leq T$). If he chooses to join the queue then - assuming throughout the discipline "first-come-first-served" - his total queueing time, from the instant of his arrival and joining to the instant of his service completion, will equal $kT + \tau + T$ (the third term being the customer's own service time). Since the decision is based on the customer's self-interest it will be considered correct if the cost of queueing does not exceed the reward. That is if, and only if, the (weak) inequality

$$R - C[(k+1)T + \tau] \geq 0 \quad (1)$$

is satisfied the newly arrived customer should join the queue.

It will be sometimes advantageous to view this from a slightly different angle. Let the occupancy or the state κ of the system at any arbitrary instant be defined as the ratio of the (future) queueing time of the last customer in the line and the service time T .

$$\kappa = \frac{kT + \tau}{T} = k + \frac{\tau}{T} \quad (2)$$

Furthermore, let a dimensionless index ν_s be defined as

$$\nu_s = \frac{R}{CT} \quad (3)$$

Inequality (1) is transformed into

$$\kappa \leq \nu_s - 1 \quad (4)$$

which is interpreted in the following manner: The incoming new customer observes the state κ of the system; if it does not exceed the value of $v_s - 1$ the customer joins the queue, otherwise the customer forgoes queueing as well as service. Now all customers act by the same criterion; hence κ can take on values in the interval $[0, v_s]$. The maximum value v_s will be realized under the following circumstances: an incoming customer encounters the system in a state $v_s - 1$, the maximum value at which the system is still accessible to new arrivals and - by the act of joining - he transforms the system into state v_s . Whenever the system is in a state within the interval $(v_s - 1, v_s]$ it is said to be inaccessible to new customers.

We note that the rule prescribing accessibility makes sense only if the inequality

$$v_s \geq 1 \quad (5)$$

pertains. In physical terms this means that the reward to be collected by the customer at the completion of his service must not fall short of the cost of time spent in service. If inequality (5) does not apply the proper policy is to refuse access to all customers and (possibly) disband the service station altogether.

The decision mode associated with the individual customer's self-interest has then given rise to a queueing model with Poisson arrivals, constant service times and finite queueing capacity (i.e. limited waiting room). Now the present model is different from those that have appeared in the literature on queueing theory in the following respect: In the typical model where finite queueing capacity makes its appearance the number of potentially available waiting spaces (the "size" of the waiting room) is assumed to be integral; occupancy too is considered to be an integer and changes in jumps whenever a customer departs or joins the queue. In the present model the service process changes (decreases) the occupancy continuously and uniformly, the total capacity is a real positive number (and not necessarily an integer) and arrivals - followed by absorption into the queue - bring about discrete changes. It is not difficult to verify that the specialization of capacity values to integers

(without imposing any further conditions) suffices to generate what was named before the "typical model with finite queueing capacity". We observe that, if the assumption of exponentially distributed service times replaces the constant service times postulate, the "typical model" cannot be transformed into a generalized model; both capacity and actual occupancy are, of necessity, integers in this case.

Occupancy is essentially equivalent to the concept of virtual waiting time (or occupation time as it is termed sometimes) introduced by Takacs; indeed, it is the ratio of occupation time to service time.

The mode of decision seeking social optimality will give rise to a queueing model of identical structure (though with one different parameter). Of course, unlike the self-optimizer, the social optimizer is not faced with a decision problem under conditions of certainty. Indeed he will have to take into account a somewhat probabilistic future, to wit, the Poisson stream of customers who will arrive at the service station. Now by the very quality of the homogeneous Poisson process the total useful information is contained in the knowledge of the arrival intensity (a parameter not relevant for the self-optimizer's decision). Hence the social optimizer, too, will at the instant of a customer's arrival, exercise control by observing the occupancy κ and make the new arrival join if, and only if, the criterion

$$\kappa \leq \nu_0 - 1 \quad (6)$$

is satisfied where ν_0 is a function of both ν_s and the traffic intensity. Adherence to such a rule will generate a system possessing a structure identical with that discussed before: Poisson arrivals, constant service times, finite and nonintegral queueing capacity.

It appears then worthwhile to delve deeper into the analysis of such a system. This will be done in the sections which follow.

IV. Some Basic Relations

We have then: a Poisson stream of incoming customers possessing arrival intensity λ ; a single service station; each customer requires exactly T time units for the completion of his service; there is limited waiting room and the occupancy κ can never exceed the constant* $\nu \geq 1$; access to the waiting line is granted to a new customer only if the occupancy does not exceed $\nu - 1$.

The traffic intensity ρ is defined as

$$\rho = \lambda T \quad (7)$$

We note that under the present model assumptions it is not necessary to put restrictions on the permissible values of ρ in order to obtain steady state conditions.

The state of the system at an arbitrary instant is specified either by the occupancy κ or by the pair (i, t) where i is the number of customers in the queue (i.e. inclusive of the customer in service) and t is the time which has already been devoted to the customer in service

$$i \begin{cases} = \kappa + \frac{t}{T} = \kappa + 1 - \frac{\tau}{T} & \text{if } \kappa > 0 \\ = 0 & \text{if } \kappa = 0 \end{cases} \quad (3)$$

$$t \begin{cases} = T - \tau & \text{if } \kappa > 0 \\ = 0 & \text{if } \kappa = 0 \end{cases} \quad (9)$$

* For the specific purposes of optimization - to be discussed in a later section - the constant ν will be assigned a subscript, e.g. ν_s or ν_o .

We are concerned with the steady state regime of our system. Let p_0 be defined as the probability of the service station being idle whereas $p_i(t)$ ($1 \leq i \leq [v] + 1$, $0 \leq t \leq T$) represents the probability density* pertaining to the elapsed service time t and the queue i .

Now consider the density $p_i(t)$ and, in particular, the change that is taking place $\Delta p_i(t)$ during a very small time interval Δt . Such change is associated with the difference of jump probabilities into, and out of, the state (i, t) , that is $p_{i-1}(t)\lambda\Delta t - p_i(t)\lambda\Delta t$.

We define

$$n = [v] \quad (10)$$

$$\theta = T(v-n) \quad (11)$$

take cognizance of the feasible values of i and of θ , and apply the idea of associating the density change within a small time duration with the difference of jump probabilities. The set of differential equations, pertaining to the present queueing model, is derived

$$\frac{dp_i(t)}{dt} = -\lambda p_i(t) \quad (0 \leq t \leq T) \quad (12a)$$

$$\frac{dp_i(t)}{dt} = \lambda[p_{i-1}(t) - p_i(t)] \quad (0 \leq t \leq T, 1 \leq i \leq n) \quad (12b)$$

$$\frac{dp_n(t)}{dt} = \lambda p_{n-1}(t) \quad (0 \leq t \leq T-\theta) \quad (12c)$$

$$\frac{dp_n(t)}{dt} = \lambda[p_{n-1}(t) - p_n(t)] \quad (T-\theta \leq t \leq T) \quad (12d)$$

$$\frac{dp_{n+1}}{dt} = \lambda p_n(t) \quad (T-\theta \leq t \leq T) \quad (12e)$$

* This is, of course, a joint density - it should be noted that one random variable (elapsed service time) is continuous while the other (queue size) is discrete. The representation in such terms, $p_i(t)$, possesses some advantage - for our present purpose - over a representation by a density associated with a single random variable, e.g. $\phi(\kappa)$. Potential concentrations and discontinuities (and, in fact, there is a concentration at the point $\kappa=0$ and a discontinuity in $\phi(\kappa)$ at the point $\kappa=1$) will be exhibited in a more natural way on utilizing the present notation.

Boundary conditions are established on examination of the changes that take place at times $t = 0$ and $t = \theta$

$$p_0 \lambda = p_1(T) \quad (13a)$$

$$p_1(0) = p_0 \lambda + p_2(T) \quad (13b)$$

$$p_i(0) = p_{i+1}(T) \quad (1 < i \leq n) \quad (13c)$$

$$p_{n+1}(T-\theta) = 0 \quad (13d)$$

The probability of having i (> 0) customers in the queue is given by

$$p_i = \int_0^T p_i(t) dt \quad 1 \leq i \leq n \quad (14a)$$

$$p_{n+1} = \int_{T-\theta}^T p_{n+1}(t) dt \quad (14b)$$

Obviously these probabilities - together with p_0 - obey

$$\sum_{i=0}^{n+1} p_i = 1 \quad (15)$$

The probability, p_c , of the service station being closed to incoming traffic may be evaluated as

$$p_c = \int_0^{T-\theta} p_n(t) dt + \int_{T-\theta}^T p_{n+1}(t) dt \quad (16)$$

The busy fraction - which in the present type of model is not identical with the traffic intensity ρ - is equal to

$$b = 1 - p_0 \quad (17)$$

During the busy period the rate of discharge of customers from the service station equals T^{-1} . Hence the average rate of discharge - i.e. the expected

number of customers leaving the service station in unit time is then evaluated as the product $(1-p_c)T^{-1}$. Now the average number of customers admitted to the service station in unit time is given by $\lambda(1-p_c)$. Within a steady state regime these two quantities must be equal. Hence after some rearrangement we obtain

$$\rho = \lambda T = \frac{1 - p_0}{1 - p_c} \quad (18)$$

If, as is assumed here, service times are fixed and equal then, by first principles, the average number of times a state t (disregarding i) is realized in unit time cannot depend on t . Hence the solutions $p_i(t)$ must obey the following equation

$$\sum_{i=1}^{n \text{ or } n+1} p_i(t) = (1-p_0)T^{-1} = \lambda(1-p_c) \quad (19)$$

It is apparent that the idle fraction p_0 plays an important role in the central formulas of the model. This quantity is a function of the parameters v and ρ . Depending on the circumstances we may desire to use the obvious notation $p_0(v, \rho)$ or $p_0(v)$.

Application of (rather lengthy) standard solution methods on the set {12} as well as combination with other equations of this Section yields

$$\begin{aligned} p_0(v, \rho) &= \left\{ 1 + \sum_{j=1}^n (-1)^{j-1} \frac{[(n-j)\lambda T + \lambda \theta]^{j-1}}{(j-1)!} \lambda T^{(n-j)\lambda T + \lambda \theta} \right\}^{-1} = \\ &= \left\{ 1 + \sum_{j=1}^n (-1)^{j-1} \frac{[(v-j)\rho]^{j-1}}{(j-1)!} \rho e^{(v-j)\rho} \right\}^{-1} \end{aligned} \quad (20)$$

If n exceeds the value 1 (the alternative case is elementary) it may be shown after some further manipulation that the following result is attained

$$p_1(v, \rho) = p_0(v, \rho)(e^\rho - 1) \quad \text{if } n > 1 \quad (21a)$$

$$p_i(v, \rho) = p_0(v, \rho) \left[e^{i\rho} + \sum_{j=1}^{i-1} (-1)^j e^{(i-j)\rho} \left\{ \frac{[(i-j)\rho]^j}{j!} + \frac{[(i-j)\rho]^{j-1}}{(j-1)!} \right\} \right] \quad \text{for } 1 < i < n \quad (21b)$$

Set (21) is rather interesting; formally the probabilities are given by equations which are identical with those relating to the analogous model with unlimited waiting room. These were already evaluated in the early days of queueing theory - indeed they can be found in Fry's (1928) textbook. However, beyond the formal identity we must take note that the two positions diverge in three aspects at least: a) The probability p_0 which appears as a multiplier in (21) is different in the two cases. b) The traffic intensity ρ must fall short of the value 1 in the infinite waiting room model; in the limited waiting room model this restriction is removed. c) In the infinite waiting room model the validity of (21b) ranges over all feasible values of i ; in the present case the range of applicability of (21b) is limited to those values of i for which the station is never closed. Beyond the intrinsic usefulness of the set {21} we are made to realize - through its presentation - that basic formulas may be stable in some sense even though some model assumptions are modified - slightly or otherwise. The change brought about by the model modification manifests itself only in the variation of a key quantity, e.g. in the present case: p_0 .

It may be desirable to examine the solution of the set {12} of differential equations in somewhat sharper detail. The following is a representation of $p_i(t)$.

Let two functions $z_i(t)$ and $Z_i(t)$ be recursively defined as

$$z_i(t) = e^{\lambda T} z_{i-1}(0) + \lambda [Z_{i-1}(t) - Z_{i-1}(T)] \quad i > 2 \quad (22)$$

$$Z_i(t) = \int_0^t z_i(t') dt' \quad i \geq 2 \quad (23)$$

and let the "starting function $z_2(t)$ " be equal to

$$z_2(t) = e^{\lambda T} - [\lambda(T-t) + 1] \quad (24)$$

The solution $p_i(t)$ of the set {12} is then given by

$$p_0(t) = p_0 e^{\lambda(T-t)} \quad (25a)$$

$$p_i(t) = p_0 e^{\lambda(T-t)} z_i(t) \quad (25b)$$

for all i and t in $(1 \leq i < n, 0 \leq t \leq T)$ which is feasible for $v \geq 2$ and for i and t in $(i=n, T-\theta \leq t \leq T)$ in which case v may take on any value exceeding 1. The restrictions on i, t and v enumerated above may be physically interpreted as relating the set {25} to precisely those states in which the service station is a) busy, and b) accessible. The proof of {25} is inductive and rather lengthy; it will not be presented here.

Finally, in this Section we put forward* an equation representing the average queue size, $q(v, \rho)$, in its dependence on v and ρ

$$\begin{aligned}
 q(v, \rho) &= n - p_0(v, \rho) \{ e^{(n-1)\lambda T + \lambda \theta} (1 - \lambda \theta) + \sum_{j=2}^n (e^{(n-j)\lambda T + \lambda \theta} \left[\sum_{l=1}^{j-1} (-1)^{l-1} \frac{[(n-j)\lambda T + \lambda \theta]^{l-1}}{(l-1)!} + (-1)^{j-1} \frac{[(n-j)\lambda T + \lambda \theta]^{j-1}}{(j-1)!} (1 - \lambda \theta) \right] \} \} = \\
 &= n - p_0(v, \rho) \{ e^{(v-1)\rho} (1 + n\rho - v\rho) + \sum_{j=2}^n (e^{(v-j)\rho} \left[\sum_{l=1}^{j-1} (-1)^{l-1} \frac{[(v-j)]^{l-1}}{(l-1)!} + \right. \\
 &\quad \left. + (-1)^{j-1} \frac{[(v-j)\rho]^{j-1}}{(j-1)!} (1 + n\rho - v\rho) \right] \} \} \quad (26)
 \end{aligned}$$

In equation (26) sums are defined to equal zero if the lower value of the summation index exceeds the upper one, e.g. $\sum_{j=2}^1 () = 0$. Hence the queue formula (26) is valid for all values of $v \geq 1$.

* Here again no proof is furnished in the paper; we wish to state that the derivation of (26) is burdensome and apparently manipulative skill rather than depth is required.

V. Optimization

The derivation of a strategy for self-optimization is rather elementary. The self-optimizing customer is aware of the quantities R , C and T . He utilizes relation (3) to compute v_s . Upon arrival at the service station he observes the actual occupancy κ . If inequality (4) is observed he reaches an affirmative decision to join. The decision is negative in the alternative case.

The impact of this strategy on the "society" of customers is that average gross gains ensue at the rate $RT^{-1}(1-p_0)$ in unit time; the resulting congestion incurs an average cost of Cq . Hence the average net gain, P , accruing to customers in unit time is given by

$$P = RT^{-1}(1-p_0) - Cq \quad (27)$$

The quantities p_0 and q in (27) are computed through the use of equations (20) and (26), the arguments ρ and v in these equations are the observed traffic intensity λT and the chosen strategy v_s , respectively.

Next we consider social optimization. Our point of departure is equation (27) and it is presently assumed that p_0 and q are functions of ρ (an observed datum) and of a v whose optimal value, v_0 , will have to be determined.

Now G is a continuous (and possibly unimodal) function of v and hence the technique of optimization that suggests itself is differentiation. After surveying the structure of p_0 and of q one is prone to think that, prima facie, differentiation would be a formidable task - technically speaking. In order to obviate the technical difficulties we proceed as follows:

Two quantities, N and D , are defined as

$$N = e^{(v-1)\rho} (1+n\rho-v\rho) + \sum_{j=2}^n (e^{(v-j)\rho}) \left[\sum_{k=1}^{j-1} (-1)^{k-1} \frac{[(v-j)\rho]^{k-1}}{(k-1)!} + (-1)^{j-1} \frac{[(v-j)\rho]^{j-1}}{(j-1)!} (1+n\rho-v\rho) \right] \quad (28)$$

$$D = 1 + \sum_{j=1}^n (-1)^{j-1} \frac{[(v-j)\rho]^{j-1}}{(j-1)!} \rho e^{(v-j)\rho} \quad (29)$$

Using this notation we may write

$$p_o = \frac{1}{D} \quad (30)$$

$$q = n - \frac{N}{D} \quad (31)$$

The derivatives of N and D (with respect to θ) are closely related

$$\frac{dN}{d\theta} = -\theta K(\theta, n) \quad (32)$$

and

$$\frac{dD}{d\theta} = TK(\theta, n) \quad (33)$$

where the function $K(\theta, n)$ is defined as

$$K(\theta, n) = \lambda^2 e^{(n-1)\rho + \lambda\theta} + \sum_{j=2}^n (-1)^{j-1} \lambda^2 e^{(n-j)\rho + \lambda\theta} \left(\frac{[(n-j)\rho + \lambda\theta]^{j-1}}{(j-1)!} + \frac{[(n-j)\rho + \lambda\theta]^{j-2}}{(j-2)!} \right) \quad (34)$$

Again in (34) the summation is defined to yield zero if the lower value of the index exceeds the upper one. Hence, $K(\theta, n)$ is defined over all feasible values of the arguments n and θ . It is not difficult to verify that it never takes on negative values.

We obtain the derivative, with respect to θ , of the net profit function (n is held constant, of course)

$$\begin{aligned} \frac{dG}{d\theta} &= \frac{d}{d\theta} \left[\frac{R(1-p_o)}{T} - Cq \right] = \frac{d}{d\theta} \left[\frac{R}{T} - \frac{R}{TD} - Cn + \frac{CN}{D} \right] = \\ &= \frac{K(\theta, n)}{D^2} [R - C(D\theta + NT)] = \frac{K(\theta, n)}{D} \left[\frac{R}{D} - C(\theta + T \frac{N}{D}) \right] = \\ &= \frac{K(\theta, n)}{D} [Rp_o - C[\theta + T(n-q)]] = \frac{CTK(\theta, n)}{D} (p_o v_s - v + q) \end{aligned} \quad (35)$$

The quantity

$$D(p_0 v_s - v + q) = v_s - \frac{D\theta}{T} - N \quad (36)$$

is a uniformly decreasing function of θ since its derivative with respect to θ is made to equal

$$\frac{d(v_s - \frac{D\theta}{T} - N)}{d\theta} = -\frac{D}{T} - \frac{\theta}{T} \left(\frac{dD}{d\theta} + \frac{T}{\theta} \frac{dN}{d\theta} \right) = -\frac{D}{T} \quad (37)$$

on utilizing equations (32) and (33).

Now for sufficiently small v expression (36) can be made positive (given that $v_s > 1$) and for sufficiently large v we can always make it negative. As all other factors on the right hand side of (35) are positive, we deduce that the function $\frac{dG}{dv}$ possesses exactly one zero at that value of the argument (v or θ) at which the function $v_s p_0 - v + q$ vanishes. Hence, v_0 , the value which brings about social optimality, may be obtained from

$$v_s p_0(v_0, \rho) - v_0 + q(v_0, \rho) = 0 \quad (38)$$

Equation (38) is of both theoretical and practical interest. We note that the problem was originally set in terms of obtaining the derivative (with respect to θ) of the net gain function. The differentiation would have to be carried out within strips of constant n since a change in n causes θ to jump between its two extreme feasible values, 0 and T. The analysis undertaken and, in particular, the devices utilized generated optimization equation (38) in which dependence on θ is suppressed and a simple formal structure is attained.

Formula (38) is also a convenient starting point for numerical work. As formulated in this study the determination of v_s precedes that of v_0 ; hence v_0 turns out to be an implicit (and not particularly convenient) function of v_s and ρ . To set up a table of numerical values one would start with given v_0 and ρ and seek the appropriate value of v_s . This is analogous to the approach undertaken in the previous study. The numerical computation

of v_s as a function of v_o and ρ is straightforward and presents no extraordinary practical difficulties. Furthermore, the physical interpretation of such a reformulation of (38) ($v_s = [v_o - q(v_o, \rho)] / p_o(v_o, \rho)$) is not farfetched: For a queueing model of the type described here the traffic intensity ρ and the socially optimal capacity v_o are given; it is desired to find that capacity, v_s , which self-optimizing customers will generate, if no regulation of traffic - financial or administrative - is imposed.

What is the optimal (maximal) rate of net gain G_o ? To derive this we return to (27) and assume that the optimal v , i.e. v_o , has been made the criterion of decision. We have then

$$\begin{aligned} G_o &= \frac{R(1-p_o(v_o, \rho))}{T} - Cq(v_o, \rho) = \\ &= C[v_s(1-p_o(v_o, \rho)) - q(v_o, \rho)] = \\ &= C[v_s - v_o - (v_s p_o(v_o, \rho) - v_o + q)] = C(v_s - v_o) \end{aligned} \quad (39)$$

Equation (39) is both simple and informative:

First it makes one realize in immediate terms that the inequality,

$$v_s \geq v_o \quad (40)$$

must hold (where equality is realized if, and only if, $v_s = 1$ ($\rho > 0$)). This, of course, is one of the objectives of the present study.

Secondly, we observe that the right hand side of (39) is very closely related to the regulatory toll that should be imposed on incoming customers in order to maximize average (social) net gain in unit time. Indeed, the optimal toll*, S_o , is obviously given by

$$S_o = CT(v_s - v_o) \quad (41)$$

* Unlike the case discussed in the previous study (where an optimal toll was one taken from a range of values) there is exactly one optimal toll value which maximizes (social) net gain in unit time.

We have then the interesting (and, on first sight, slightly strange) result that the optimal toll to be imposed on the customer is identical with the average optimal gain accumulating during one service period.

$$S_0 = G_0 T \quad (42)$$

Thirdly, one is induced to pose the question whether the simple formulas attained here - such as (38), (39) and (41) - are amenable to simple physical interpretations and, possibly, to further generalizations. In the following we present the marginal analysis pertaining to social optimization. We shall show that it leads to the very same equations possessing elementary structure.

It is the social optimizer's function to select an indifference capacity ($v_0 - 1$) possessing the following characteristic: A customer who arrives at an instant - $t = 0$, say - at which the system possesses the occupancy $\kappa = v_0 - 1$ will generate identical gains to society either by joining the queue or by declining to do so. Neither alternative is preferable to the other from the viewpoint of public good. We note that if the customer joins the queue the identical state $v_0 - 1$ which would instantaneously prevail were he to balk will be regenerated* in exactly T time units (with probability 1). During this time access to the service station is blocked for new customers who (possibly) arrive within that interval. Exactly one customer will be discharged from the service station during the blocked period - at time $(v_0 - n_0)T$. The queue size before and after this discharge is $n_0 + 1$ and n_0 , respectively; it is easy to verify that the average queue length is v_0 . Hence, as a result of joining, the total net benefits reaped during T amount to $R - CTv_0$. However, the decision to join at time $t = 0$ (and occupancy $\kappa = v_0 - 1$) has further implications. It will be convenient to represent them by an instantaneous expected net gain rate $g_{\text{join}}(t; v_0 - 1)$. Since the net gain $(R - CTv_0)$ during the interval $(0, T)$ has already been separated out the function $g_{\text{join}}(t; v_0 - 1)$ takes on the value 0 up to time T

$$g_{\text{join}}(t; v_0 - 1) = 0 \quad (0 \leq t \leq T) \quad (43)$$

* This property depends on the assumptions that customers arrive in a Poisson stream at the station.

Beyond T the function $g_{\text{join}}(t)$ takes a course which incorporates the presently existing queue, the accumulation of new customers, the discharge of serviced customers and the rewards gained by them. Clearly the instantaneous expected net gain rate tends to $G(v_0)$ as time t tends to infinity

$$g_{\text{join}}(t; v_0-1) \xrightarrow{t \rightarrow \infty} P(v_0) \quad (44)$$

The alternative decision to balk at time $t = 0$ brings forth another instantaneous expected net gain rate $g_{\text{balk}}(t; v_0-1)$. By virtue of the characteristics stated before, joining the queue at time $t = 0$ generates a state at time $t = T$ which is identical with the state at time $t = 0$ brought about by the balking decision. Hence the following must hold

$$g_{\text{balk}}(t; v_0-1) = g_{\text{join}}(t+T; v_0-1) \quad (45)$$

and, of course, analogously to (44) we have

$$g_{\text{balk}}(t; v_0-1) \xrightarrow{t \rightarrow \infty} G(v_0) \quad (46)$$

What is the expected accumulated financial advantage $A(t)$ at t (conveniently assumed to exceed T) of balking over joining where we disregard the terms $R - C v_0 T$ which favored joining and were separated out. Clearly $A(t)$ is given by

$$\begin{aligned} A(t) &= \int_0^t g_{\text{balk}}(t') dt' - \int_0^t g_{\text{join}}(t') dt' = \int_0^t g_{\text{balk}}(t') dt' - \\ &- \int_T^t g_{\text{join}}(t') dt' = \int_0^t g_{\text{balk}}(t') dt' - \int_0^{t-T} g_{\text{balk}}(t') dt' = \\ &= \int_{t-T}^t g_{\text{balk}}(t') dt' \end{aligned} \quad (47)$$

When t tends to infinity the integrand on the right hand side of (47) tends to the constant $G(v_0)$; hence the integral (47) - with t tending to infinity - is evaluated as

$$A(\infty) = TG(v_0) = TG_0 \quad (48)$$

The gist of marginal analysis is that under conditions of optimality, this advantage of balking over queueing (over an infinite horizon*) must equal the advantage $R - v_0 CT$ of joining over balking within the interval $(0, T)$. Therefore, we obtain

$$R - v_0 CT = TG_0 \quad (49)$$

But equation (49) is essentially identical with (39). The other general optimization formulas (38) and (41) may be easily derived from (39). Hence by using marginal analysis we have obtained the procedure for optimization without the "messy" computational technicalities. For actual numerical work it is, of course, still necessary to evaluate queue sizes and idle fractions through the use of formulas (28)-(31).

Marginal analysis has led us one step beyond the original model under investigation. The argument leading to (49) - and hence to (38) and (41) - remains essentially valid even if the assumption of constant and equal service times is modified. It is sufficient to assume that service times are distributed (rather than fixed) and that the class of distributions is characterized by the expected remaining service time of a customer being a strictly decreasing and continuous function of elapsed service time. This is a rather mild restriction. The minor modification that has to be introduced in the argumentation is that the phrase "exactly after T time units" has to be replaced by "after T time units on the average" whenever it appears. The salient point is the following: whenever a situation exists such that a marginally joining customer is made to produce average net gain in unit time during the ensuing T time units, on the average, equations (38), (39) and (41) must hold. We mention in passing that, under conditions of social optimality a customer admitted to the service station in a non-marginal fashion as it were generates net gain exceeding the average.

* We observe that the interest rate is (implicitly) assumed to equal zero; hence it does not make an appearance in the argument.

Even if service times are distributed in a manner other than that prescribed in the preceding paragraph, equations (38), (39) and (41) may remain valid - at least in some approximative sense. Thus, for instance, let it be assumed that the service times are exponentially distributed; this is the case discussed in the previous study. Clearly the expected remaining service time of a customer is a constant rather than a strictly decreasing function as postulated before. Yet if in the equations representing the idle fraction and the queue the integer n_0 is replaced by the (close) real number v_0 it can be shown that relations (38) etc. are revalidated. At the danger of being repetitious let it be restated that the analytical, as well as numerical, derivation of the optimal p_0 and q may be quite a difficult task.

VI. Conclusion

The program of this investigation was threefold:

The first objective was to show that the decision rule of self-optimizing customers operating within a framework of certainty and of equality (pertaining to R, C and T) tends to overcongest a queueing system. This is the proper meaning of inequality (40). The basic reason for the divergence between social optimization and self-optimization - as expressed in inequality (40) - is the fact that the individual customer need not consider the penalties he is (possibly) inflicting upon future customers by the very act of his joining the queue. The toll levied on a marginally joining customer could be considered to represent compensation for damage, as it were, caused by the customer to future customers.

The second objective was to establish a vantage point for further generalization. This has been attained by alternating ordinary maximization (that is: carried out by differentiation) and marginal analysis. A set of formulas, simple and comprehensive - (38), (39) and (41) - has been shown to hold under conditions more general than originally specified.

Thirdly, the stochastic queueing model with non-integral capacity has been developed and, possibly, this may be applicable in situations other than those possessing an optimization rationale. The structure and form of associated quantities - probabilities and expectations - may be quite interesting per se and some potential industrial applications indicate the necessity for further study.

The general subject area of this study possesses useful and interesting extensions. Some further investigations are under way.

R e f e r e n c e s

- | | |
|-----------------|---|
| Thornton C. Fry | Probability and its Engineering Uses.
D. Van Nostrand Co. Inc., New York (1928). |
| P. Naor | The Regulation of Queue Size by Levying Tolls.
Econometrica, <u>37</u> , 15-24 (1969). |
| R. Shinnar | Sizing of Storage Tanks for Off-Grade Material.
Industrial and Engineering Chemistry,
Process Design and Development, <u>6</u> , 263-4, (1967). |